## Московский государственный технический университет им. Н.Э. Баумана Кафедра «Системы обработки информации и управления»



# Рубежный контроль №2 по дисциплине «Методы машинного обучения»

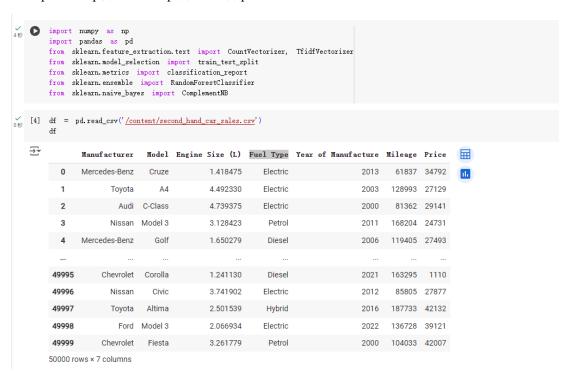
Выполнил:

студент группы ИУ5И-21М

Ван Чжэн

Москва — 2024 г.

Я выбрал набор данных-«продажа подержанных автомобилей».



## Классификатор№1: KNeighborsClassifier

KNeighborsClassifier+Tfidf

```
KNeighborsClassifier+TFIDF
🕟 from sklearn.neighbors import KNeighborsClassifier
    X = df['Manufacturer']
    y = df['Fuel Type']
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
    # 使用 TfidfVectorizer 进行向量化
    tfidf_vectorizer = TfidfVectorizer()
    X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
    X_test_tfidf = tfidf_vectorizer.transform(X_test)
    # 使用 KNeighborsClassifier 进行分类(基于 TfidfVectorizer)
    knn_tfidf = KNeighborsClassifier()
    knn_tfidf.fit(X_train_tfidf, y_train)
    y_pred_tfidf = knn_tfidf.predict(X_test_tfidf)
    print(classification_report(y_test, y_pred_tfidf, digits=4))
₹
               precision recall f1-score support
         Diesel 0.2403 0.3918
                                      0.2979
                                                  2478
                 0.2494 0.4854 0.3295
0.0000 0.0000 0.0000
        Electric
                                                 2530
         Hybrid
                                                 2547
         Petrol 0.2363 0.1002 0.1407
                                              2445
                                      0.2444
                                               10000
       accuracy
    macro avg 0.1815 0.2444 0.1920
weighted avg 0.1804 0.2444 0.1916
                                                 10000
                                                10000
```

#### KNeighborsClassifier+CountVec

#### KNeighborsClassifier+CountVec

```
  # 使用 CountVectorizer 进行向量化

    count_vectorizer = CountVectorizer()
    X_train_count = count_vectorizer.fit_transform(X_train)
    X_test_count = count_vectorizer.transform(X_test)
    knn_count = KNeighborsClassifier()
    knn_count.fit(X_train_count, y_train)
    y_pred_count = knn_count.predict(X_test_count)
    print(classification_report(y_test, y_pred_count, digits=4))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: T
      _warn_prf(average, modifier, msg_start, len(result))
                 precision
                            recall f1-score
                                               support
                    0.2360 0.3793
                                       0.2910
          Diesel
                                                   2478
        Electric 0.2540 0.4107
Hybrid 0.2425 0.1834
                                       0.3138
                                                   2530
                                       0.2088
                                                   2547
                    0.0000 0.0000
                                       0.0000
          Petrol
                                                   2445
                                       0.2446
                                                  10000
        accuracy
                    0.1831
                              0.2433
                                       0.2034
                                                  10000
       macro avg
                    0.1845 0.2446
                                       0.2047
                                                  10000
    weighted avg
```

## Классификатор№2: LogisticRegression

LogisticRegression+TFIDF

```
from sklearn.linear_model import LogisticRegression
    # 使用 LogisticRegression 进行分类(基于 TfidfVectorizer)
    lr_tfidf = LogisticRegression()
    lr_tfidf.fit(X_train_tfidf, y_train)
    y_pred_lr_tfidf = lr_tfidf.predict(X_test_tfidf)
    print(classification_report(y_test, y_pred_lr_tfidf,
                                                        digits=4))
₹
                             recall f1-score
                 precision
                                               support
                             0.3067
          Diesel
                    0.2497
                                       0.2753
                                                  2478
        Electric
                    0.2460
                             0.0976
                                       0.1398
                                                  2530
          Hybrid
                    0.2419
                             0.0876
                                       0.1286
                                                  2547
          Petrol
                    0.2507
                             0.5157
                                       0.3374
                                                  2445
                                       0.2491
                                                 10000
        accuracy
                                                 10000
                    0.2471
                             0.2519
                                       0.2203
       macro avg
    weighted avg
                    0.2470
                             0.2491
                                       0.2188
                                                 10000
```

### Logistic Regression + Count Vec

LogisticRegression+CountVec

```
▶ # 使用 LogisticRegression 进行分类(基于 CountVectorizer)
    lr_count = LogisticRegression()
    lr_count.fit(X_train_count, y_train)
    y_pred_lr_count = lr_count.predict(X_test_count)
    print(~\nLogisticRegression 基于 CountVectorizer 的分类报告: ~)
    print(classification_report(y_test, y_pred_lr_count, digits=4))
₹
    LogisticRegression 基于 CountVectorizer 的分类报告:
                precision recall f1-score
                   0.2497 0.3067
                                    0.2753
                                               2478
         Diesel
                 0.2460 0.0976
       Electric
                                    0.1398
                                               2530
         Hybrid
                  0.2419 0.0876
                                    0.1286
                                               2547
         Petrol
                 0.2507 0.5157
                                    0.3374
                                               2445
                                    0.2491
                                              10000
       accuracy
       macro avg
                 0.2471
                           0.2519
                                    0.2203
                                              10000
    weighted avg
                 0.2470
                           0.2491
                                    0.2188
                                              10000
```